# FUNDAMENTALS OF DEEP LEARNING FOR MULTI-GPUS

## LAB 1, PART 1: INTRODUCTION AND MOTIVATION

**NVIDIA** | DEEP LEARNING INSTITUTE

# COURSE OVERVIEW

- Lab 1: Gradient Descent vs Stochastic Gradient Descent, and the Effects of Batch Size

- Lab 2: Multi-GPU DL Training Implementation using Horovod

- Lab 3: Algorithmic Concerns for Training at Scale
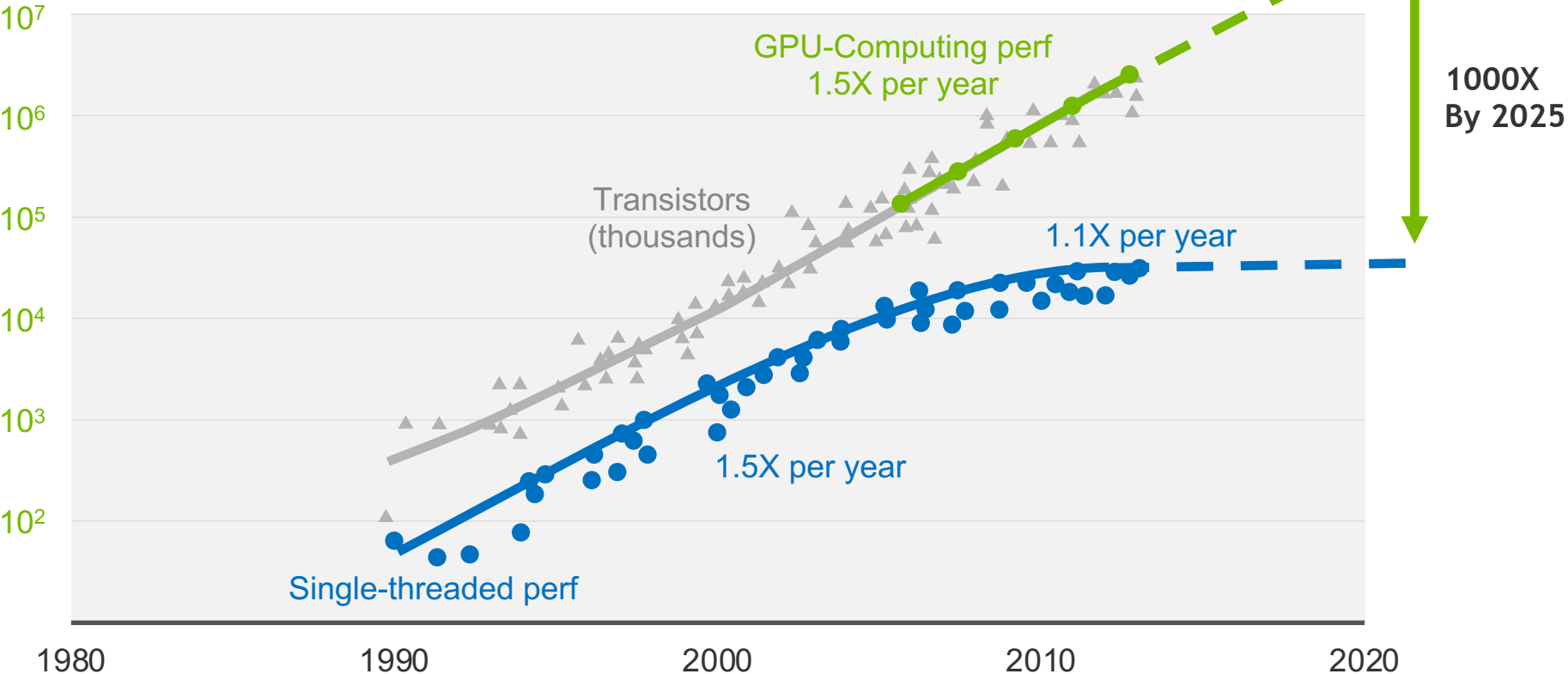
# LAB 1 OVERVIEW

- Part 1: Gradient Descent

- Part 2: Stochastic Gradient Descent

- Part 3: Optimizing training with batch size

# CONTEXT: WHY USE MULTIPLE GPUS?

# TRENDS IN COMPUTATIONAL POWER

## Historically we never had large datasets or compute



- GPU-Computing perf 1.5X per year
- Transistors (thousands)
- 1.1X per year
- 1.5X per year
- Single-threaded perf
- 1000X By 2025

$10^7$ $10^6$ $10^5$ $10^4$ $10^3$ $10^2$

1980    1990    2000    2010    2020

# TRENDS IN COMPUTATIONAL POWER
## 2 PF/s in November 2009

# TRENDS IN COMPUTATIONAL POWER
## 5 PF/s today



10x NVIDIA® Mellanox® ConnectX-6
200 Gb/s Network Interface

500 GB/s Peak Bi-directional Bandwidth

Dual 64-Core AMD Rome CPUs
2 TB RAM

3.2X More Cores to Power the Most Intensive AI Jobs

8x NVIDIA A100 Tensor Core GPUs

Up to 640 GB Total GPU Memory
12 NVIDIA NVLinks™ per GPU
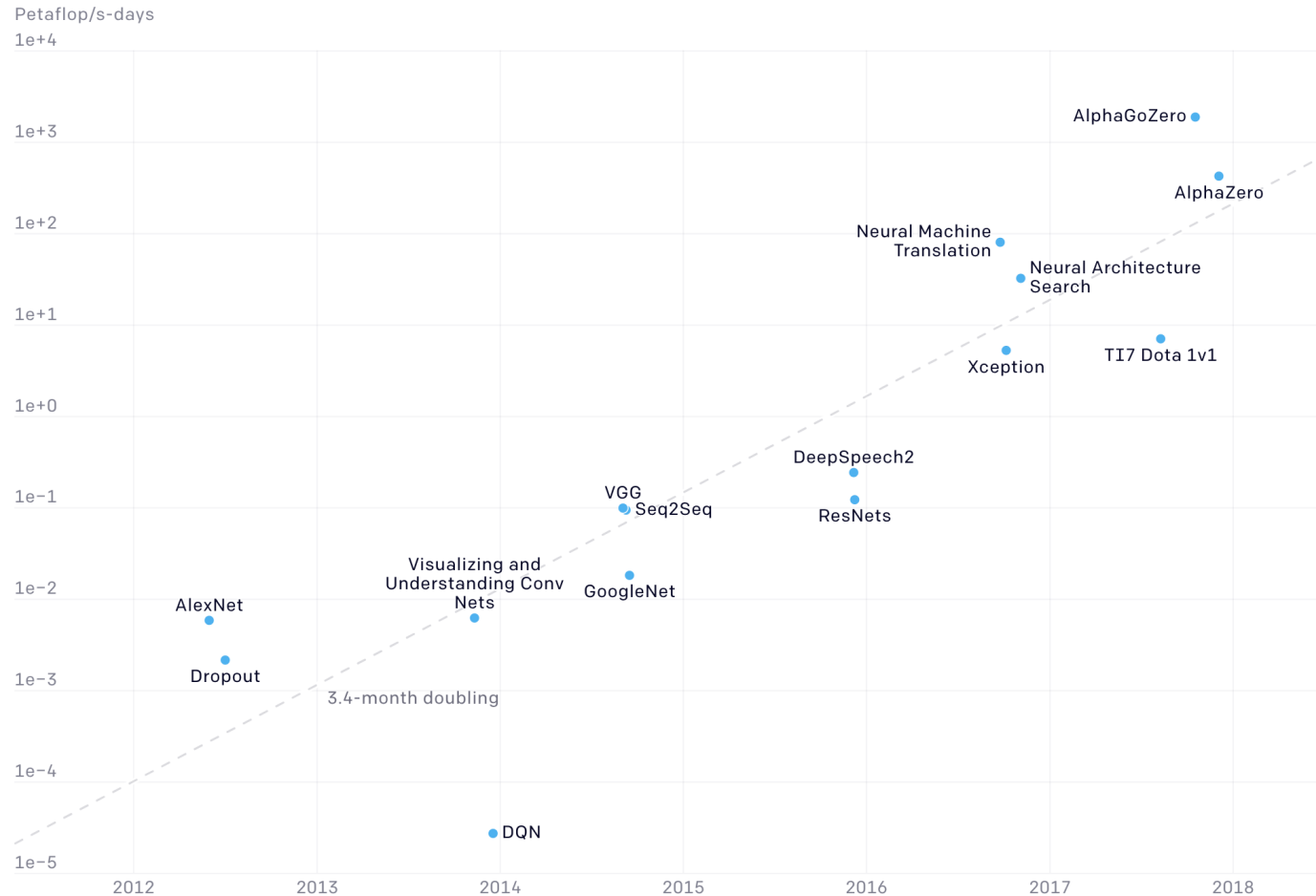600 GB/s GPU-to-GPU Bi-directional Bandwidth

6x NVIDIA NVSwitches™

4.8 TB/s Bi-directional Bandwidth
2X More than Previous-Generation NVSwitch

30 TB Gen4 NVME SSD

50 GB/s Peak Bandwidth
2X Faster than Gen3 NVME SSDs

| 1 TRILLION Transistors | 1 KILOMETER of Traces | 1 MILLION Drill Holes | 30 THOUSAND Components | 50 THOUSAND Connector Pins | 5 PETAFLOPS of AI Performance |
| --- | --- | --- | --- | --- | --- |

# NEURAL NETWORK COMPLEXITY IS EXPLODING

**AlexNet to AlphaGo Zero: A 300,000x Increase in Compute (Log Scale)**
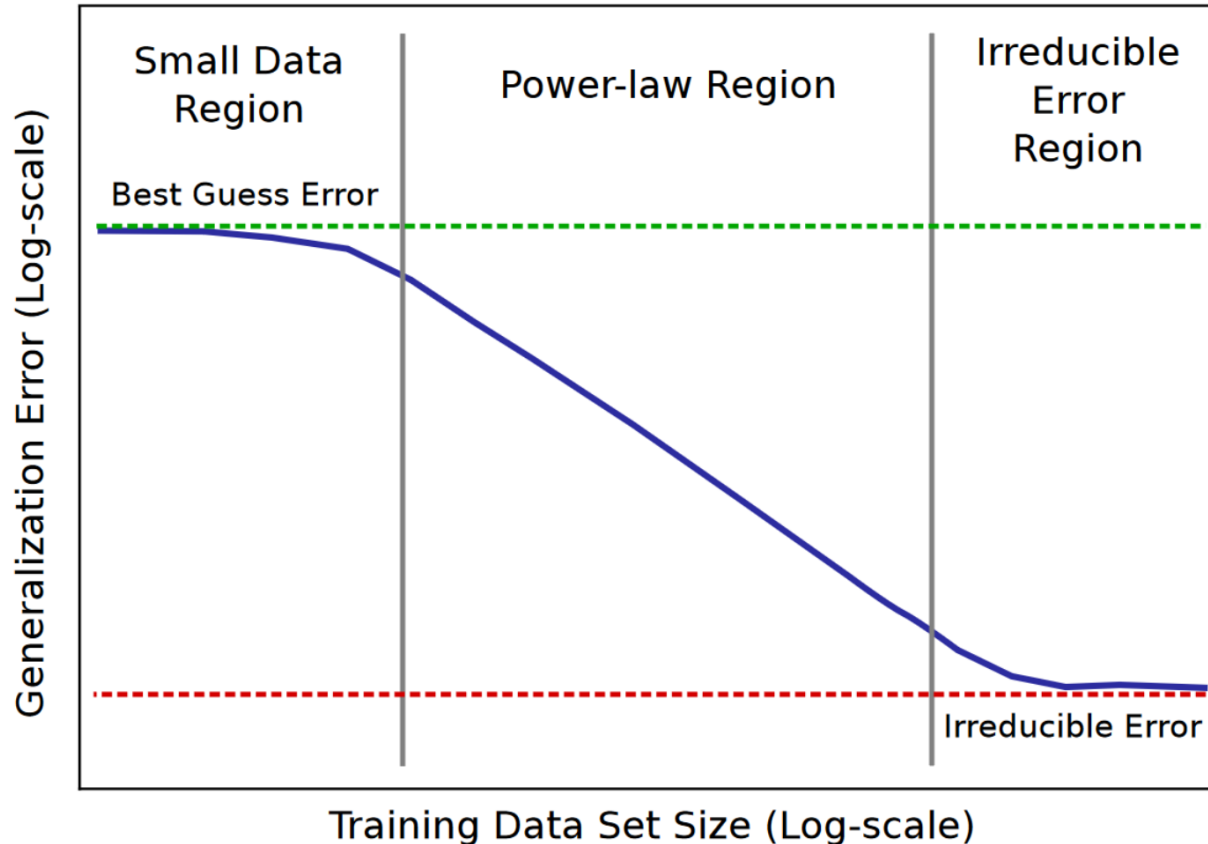
Petaflop/s-days



3.4-month doubling

Source: OpenAI

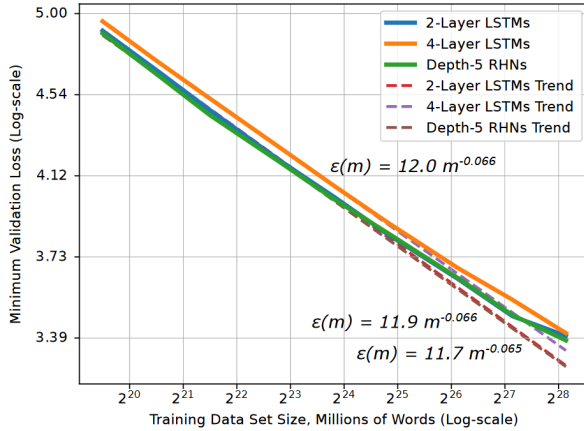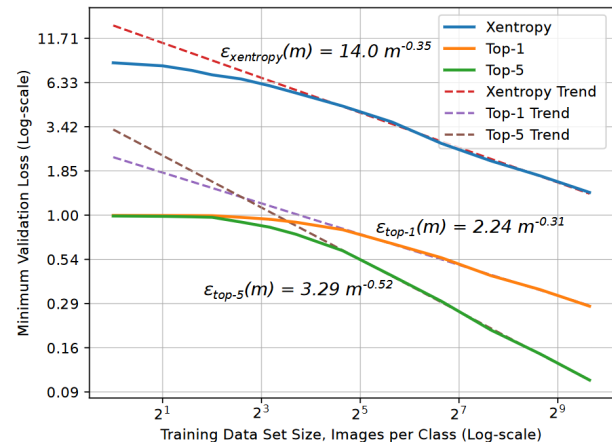# 1000 PETAFLOP/S-DAYS
# =
# O(100 YEARS) ON A DUAL CPU SERVER

# EXPLODING DATASETS
## Power-law relationship between dataset size and accuracy



Hestness, J., et al. (2017). Deep Learning Scaling is Predictable, Empirically. arXiv: 1712.00409
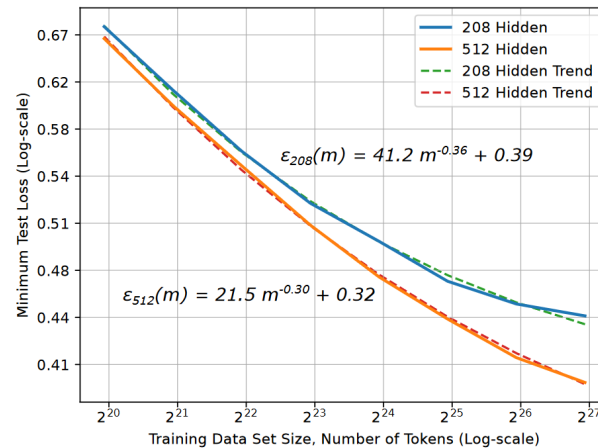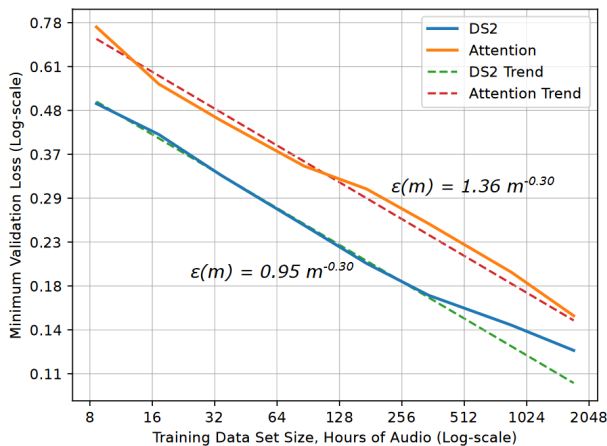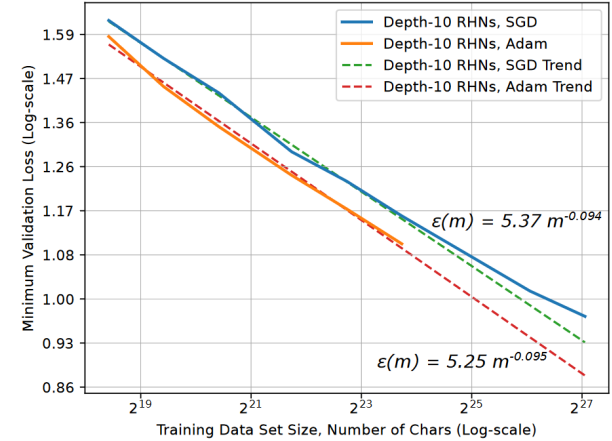
# EXPLODING DATASETS

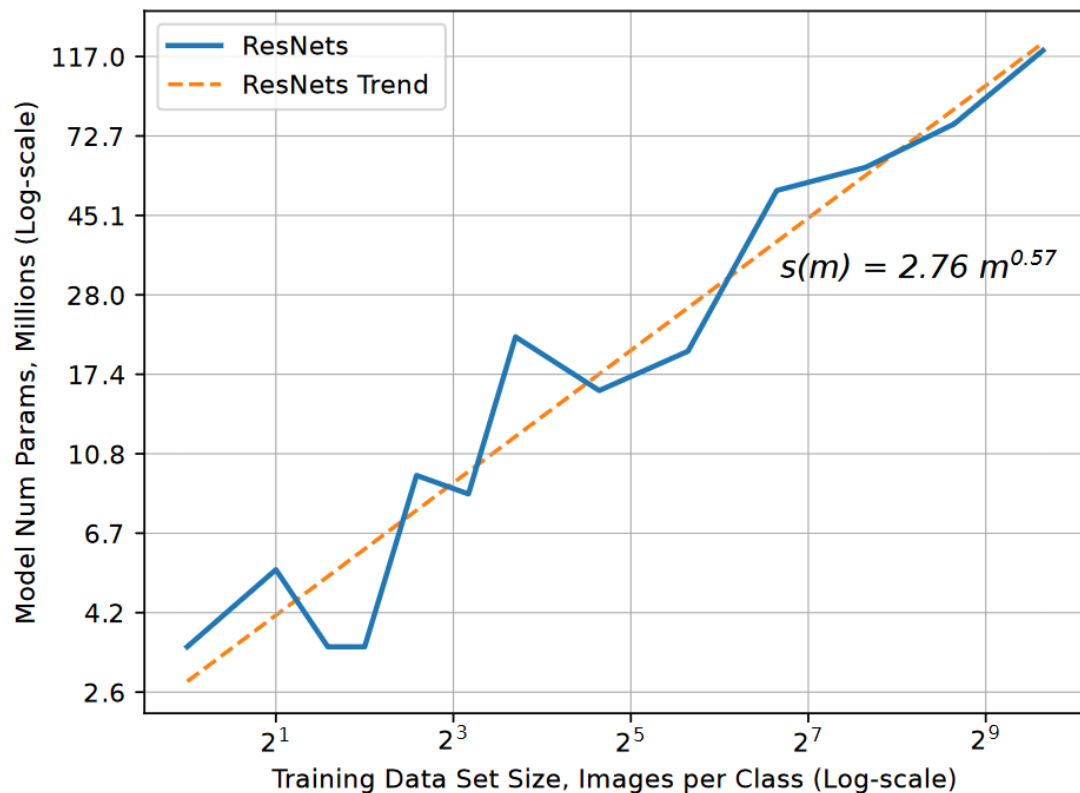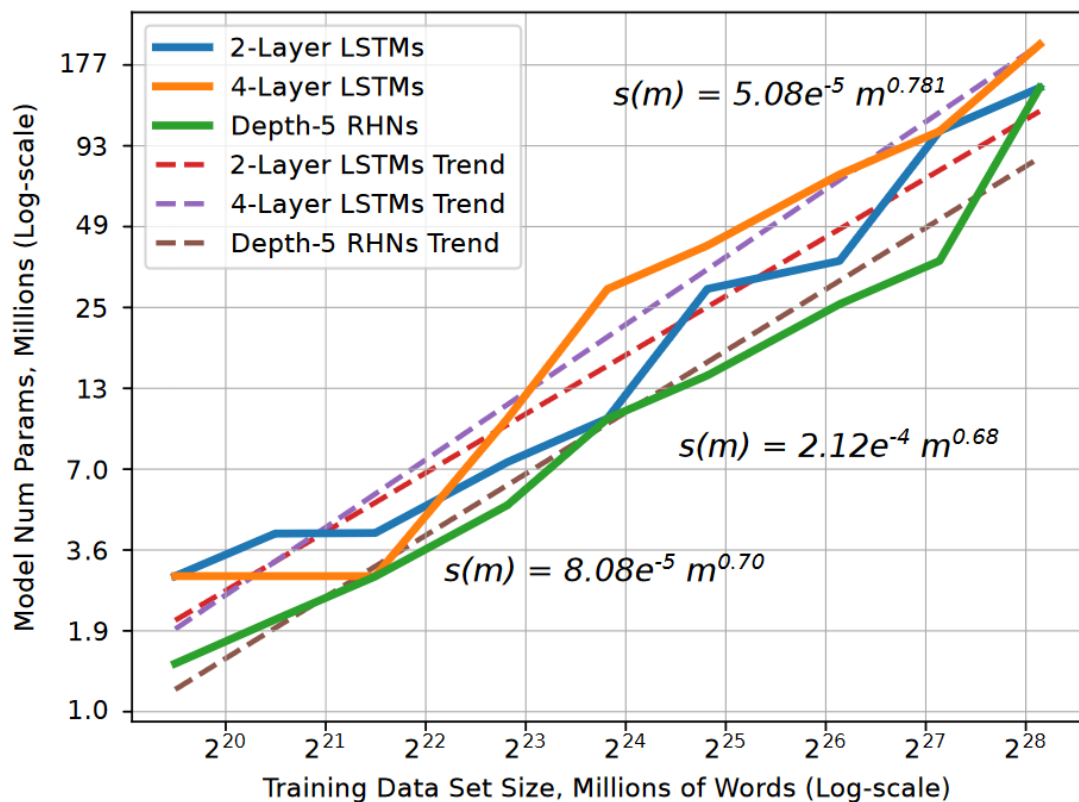## Power-law relationship between dataset size and accuracy



- Translation
- Language Models
- Character Language Models
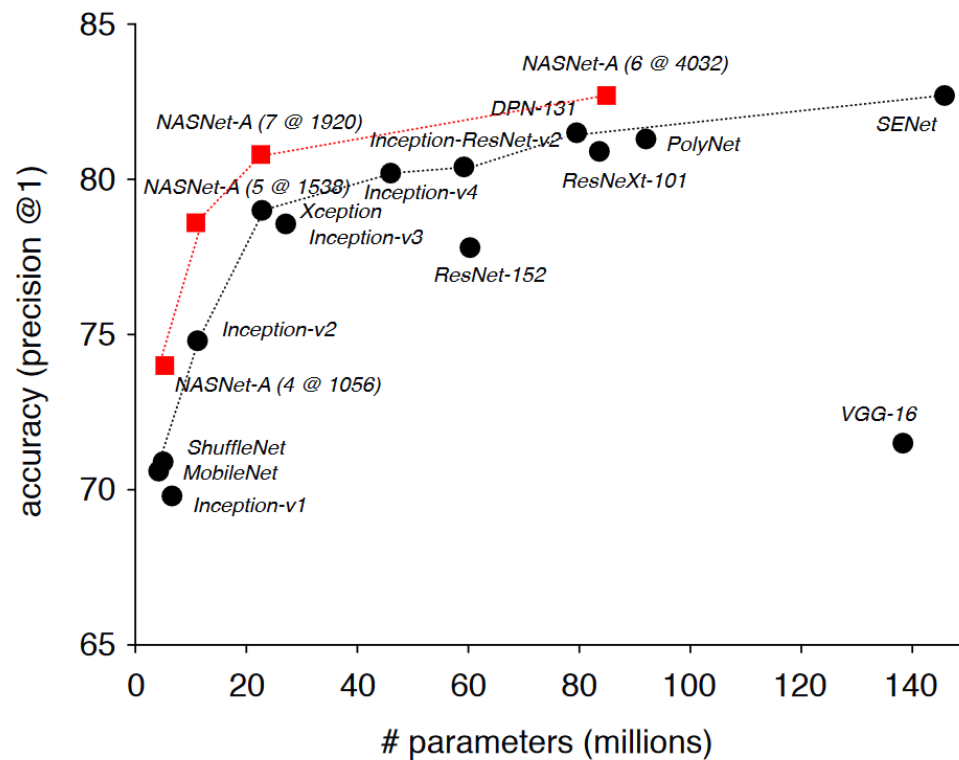- Image Classification
- Attention Speech Models

Hestness, J., et al. (2017). Deep Learning Scaling is Predictable, Empirically. arXiv: 1712.00409

# EXPLODING MODEL COMPLEXITY

## Though model size scales sublinearly



$s(m) = 5.08e^{-5} m^{0.781}$

$s(m) = 2.12e^{-4} m^{0.68}$

$s(m) = 8.08e^{-5} m^{0.70}$

$s(m) = 2.76 m^{0.57}$

Legend (left plot):
- 2-Layer LSTMs
- 4-Layer LSTMs
- Depth-5 RHNs
- 2-Layer LSTMs Trend
- 4-Layer LSTMs Trend
- Depth-5 RHNs Trend

Left plot axes: Model Num Params, Millions (Log-scale) vs Training Data Set Size, Millions of Words (Log-scale)

Legend (right plot):
- ResNets
- ResNets Trend

Right plot axes: Model Num Params, Millions (Log-scale) vs Training Data Set Size, Images per Class (Log-scale)

Hestness, J., et al. (2017). Deep Learning Scaling is Predictable, Empirically. arXiv: 1712.00409

# EXPLODING MODEL COMPLEXITY
## Though model size scales sublinearly



Zoph, Barret, et al. (2017). "Learning transferable architectures for scalable image recognition." arXiv: 1707.07012

# IMPLICATIONS
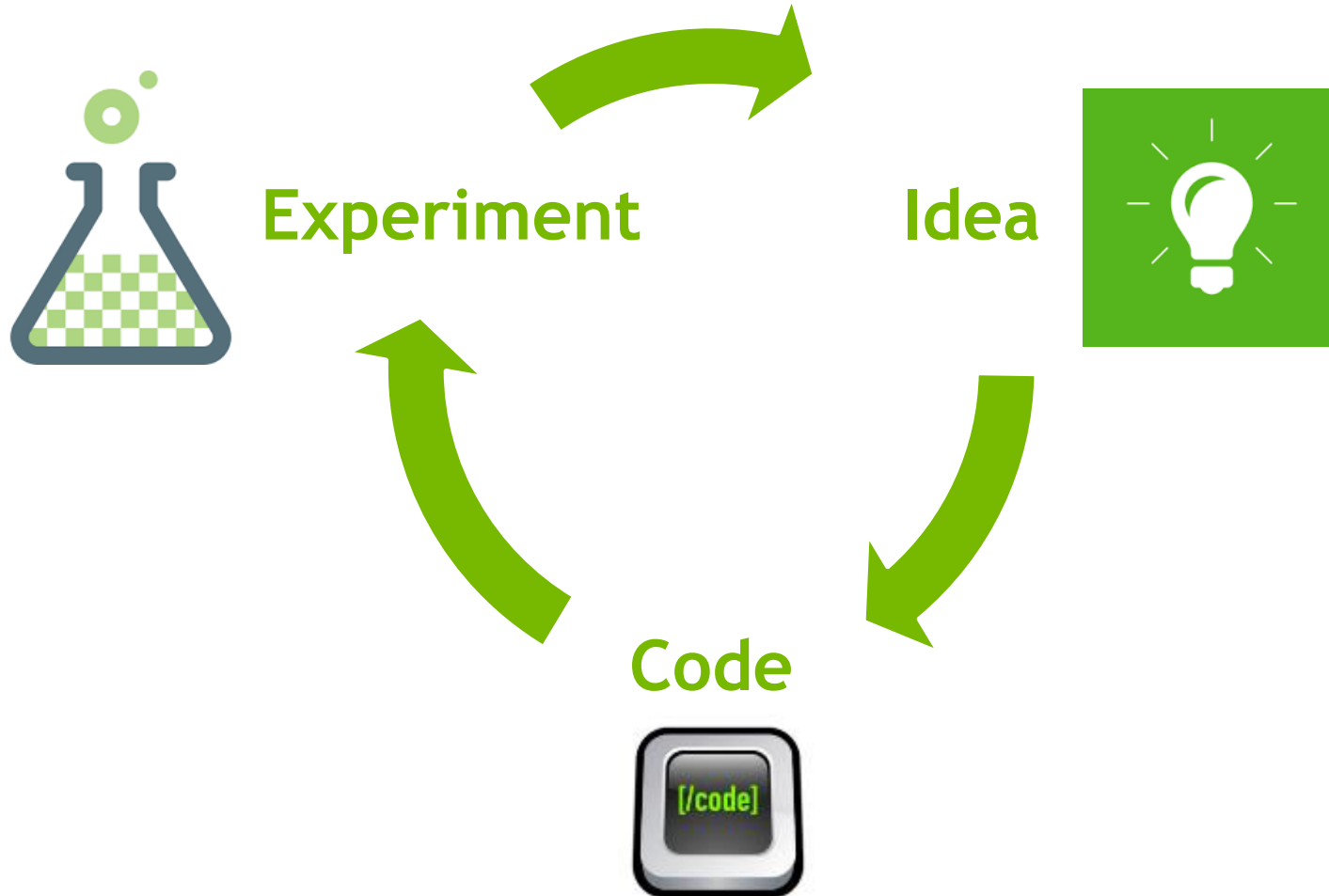
# IMPLICATIONS
## Good and bad news

▸ The good news: Requirements are predictable.

    ▸ We can predict how much data we will need.

    ▸ We can predict how much computing power we will need.

▸ The bad news: The values can be significant.

    ▸ The silver lining is that deep learning has taken impossible problems and made them merely expensive.
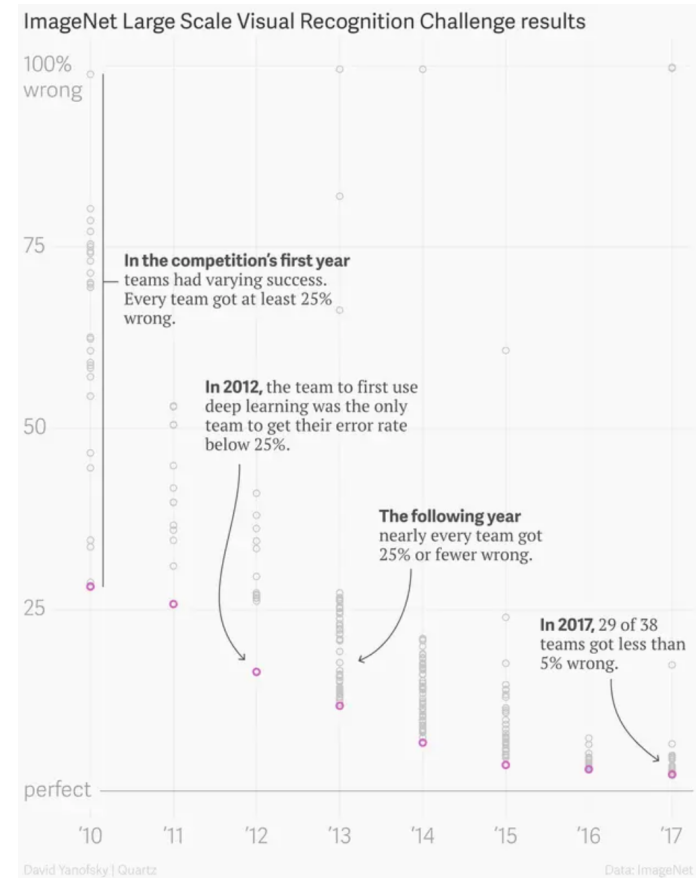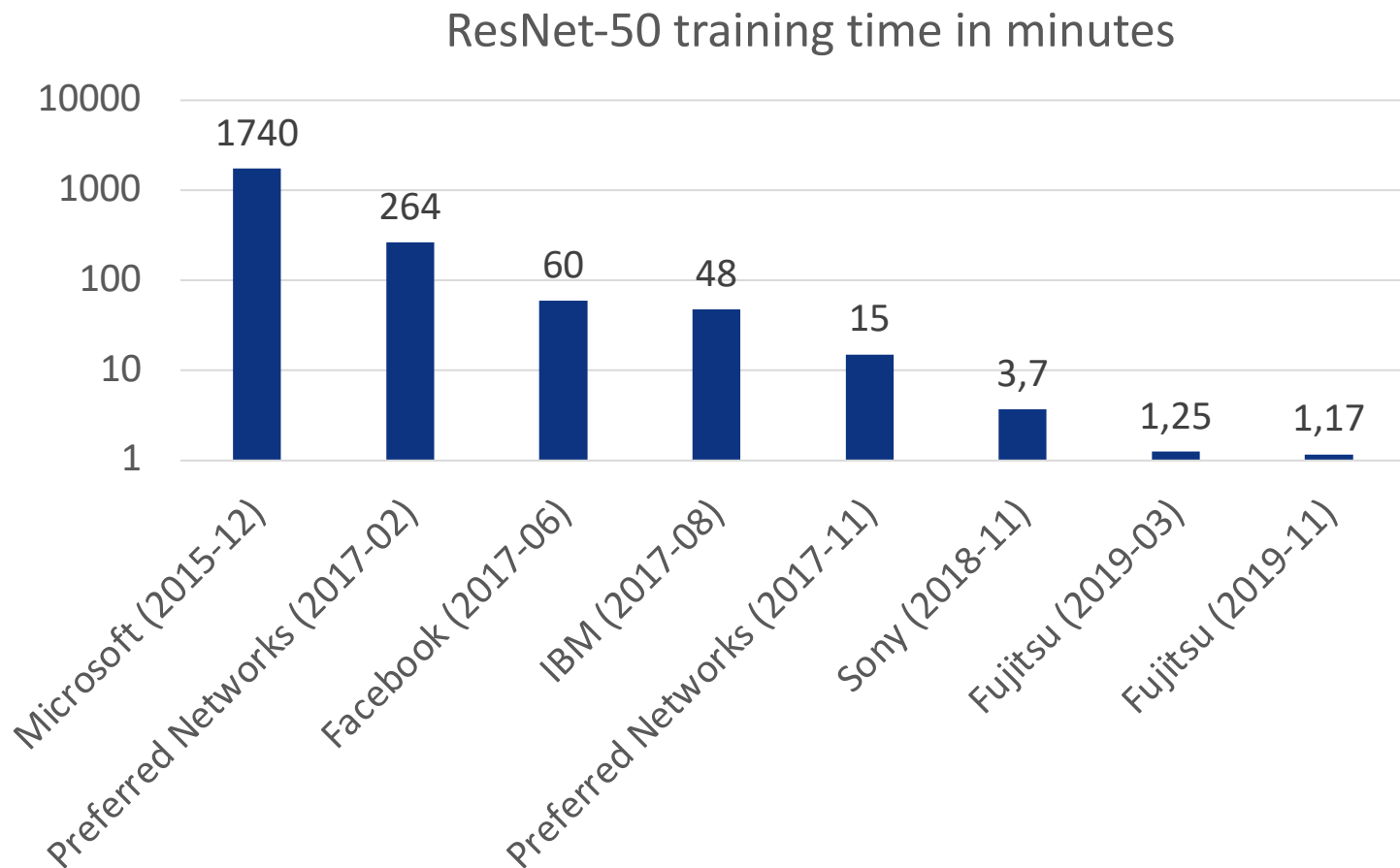
# IMPLICATIONS

Deep learning is experimental; we need to train quickly to iterate



Experiment

Idea

Code

# ITERATION TIME
## Short iteration time is fundamental for success

ResNet-50 training time in minutes

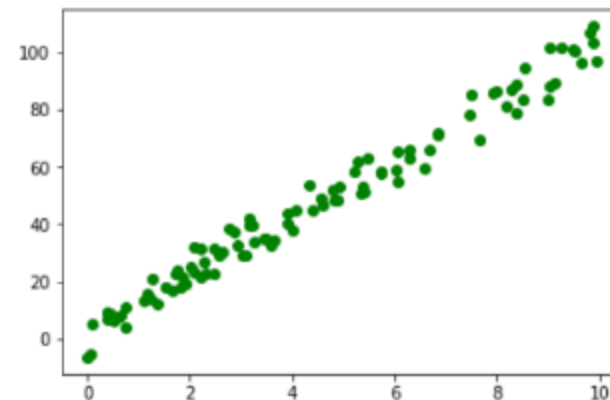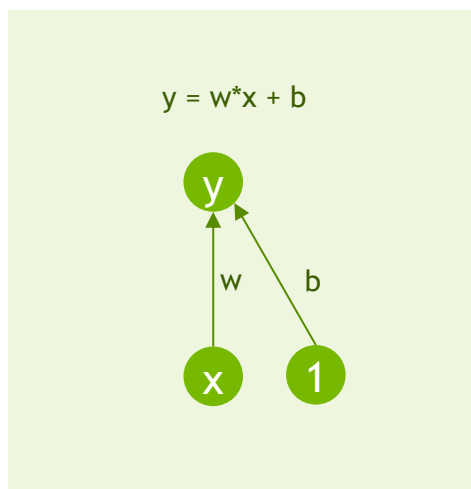

ImageNet Large Scale Visual Recognition Challenge results

# INTRO TO THE LAB

# STARTING WITH A LINEAR MODEL

Our goal is to find best model parameters (combination of w and b) to fit the data

# FUNDAMENTALS OF DEEP LEARNING FOR MULTI-GPUS

## LAB 1, PART 2: MORE REALISTIC NETWORKS
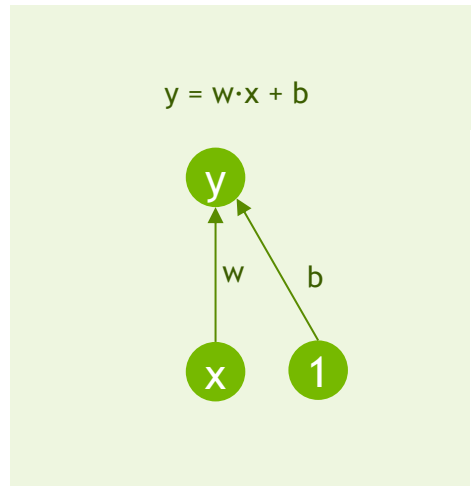
NVIDIA
DEEP
LEARNING
INSTITUTE

# MODERN NEURAL NETWORKS
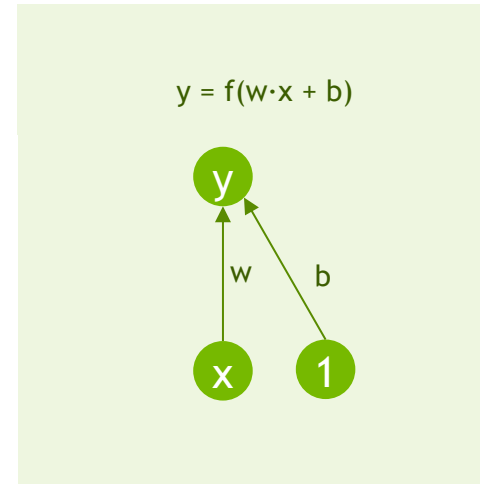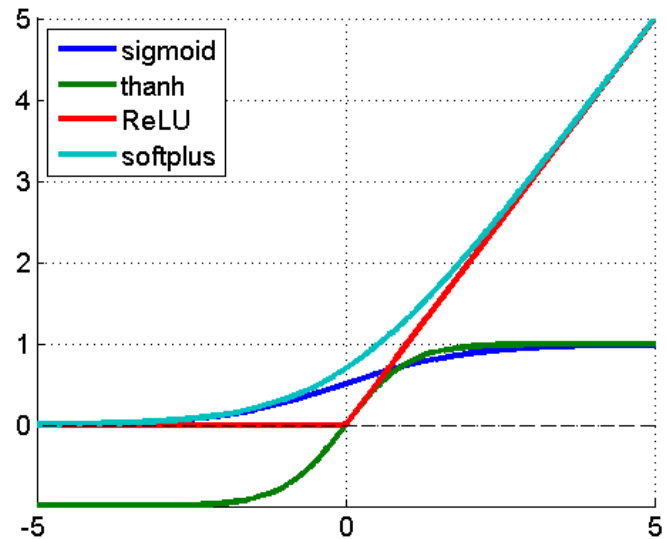
## How do they differ from our trivial example?

Not significantly!

# MODERN NEURAL NETWORKS
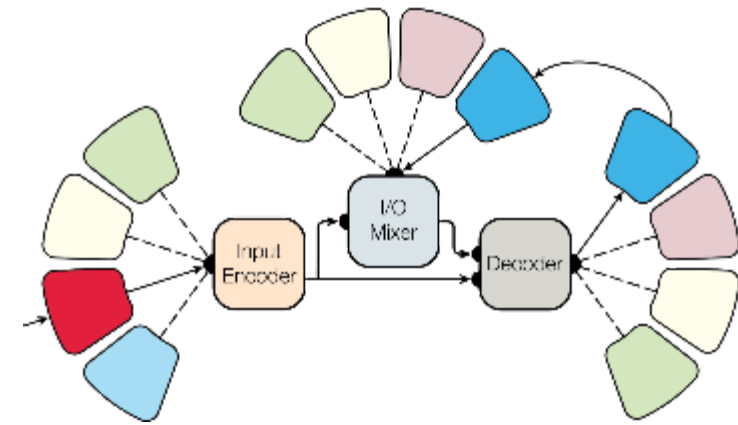
## How do they differ from our trivial example?
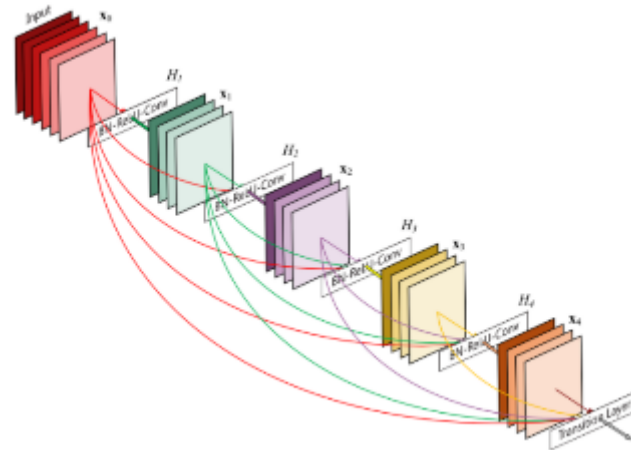
Nonlinearity

y = w·x + b

y = f(w·x + b)

# MODERN NEURAL NETWORKS

## How do they differ from our trivial example?

More complex interconnection and many more parameters

Kaiser, L., Gomez, A. N., Shazeer, N., Vaswani, A., Parmar, N., Jones, L., & Uszkoreit, J. (2017). One model to learn them all. *arXiv preprint arXiv:1706.05137.*
Iandola, F., Moskewicz, M., Karayev, S., Girshick, R., Darrell, T., & Keutzer, K. (2014). Densenet: Implementing efficient convnet descriptor pyramids. arXiv preprint arXiv:1404.1869.
Shazeer, N., Mirhoseini, A., Maziarz, K., Davis, A., Le, Q., Hinton, G., & Dean, J. (2017). Outrageously large neural networks: The sparsely-gated mixture-of-experts layer. *arXiv preprint arXiv:1701.06538.*

# NON-CONVEX LOSS FUNCTIONS
## Those differences make the optimization problem much more difficult

# NON-CONVEX LOSS FUNCTIONS
## Those differences make the optimization problem much more difficult

Linear model loss function

ResNet-56 loss function projection to 3D – no skip connections



Li, H., Xu, Z., Taylor, G., & Goldstein, T. (2017). Visualizing the Loss Landscape of Neural Nets. *arXiv:1712.09913*.

# NON-CONVEX LOSS FUNCTIONS
## Those differences make the optimization problem much more difficult

ResNet-56 loss function projection to 3D – no skip connections



Why do we succeed in finding good local minima?

Li, H., Xu, Z., Taylor, G., & Goldstein, T. (2017). Visualizing the Loss Landscape of Neural Nets. *arXiv:1712.09913*.

# NON-CONVEX LOSS FUNCTIONS

## Recent advances such as residual connections simplify optimization



(a) without skip connections

(b) with skip connections

Li, H., Xu, Z., Taylor, G., & Goldstein, T. (2017). Visualizing the Loss
Landscape of Neural Nets. *arXiv:1712.09913*.

# FUNDAMENTALS OF DEEP LEARNING FOR MULTI-GPUS

## LAB 1 CONCLUSION: DATA AND MODEL PARALLELISM

NVIDIA
DEEP LEARNING INSTITUTE

# DATA PARALLELISM
## Focus of this course

How can we take advantage of multiple GPUs to reduce the training time?

# DATA VS MODEL PARALLELISM
## Comparison

▸ Data Parallelism

- ▸ Allows you to speed up training

- ▸ All workers train on different data

- ▸ All workers have the same copy of the model

- ▸ Neural network gradients (weight changes) are exchanged

▸ Model Parallelism

- ▸ Allows you to use a bigger model

- ▸ All workers train on the same data

- ▸ Parts of the model are distributed across GPUs

- ▸ Neural network activations are exchanged

# TRAINING A NEURAL NETWORK

## Single GPU

GPU

$\hat{y} \longrightarrow \mathcal{L}(\hat{y}, y)$

$W^{[3]}$

$W^{[3]} = W^{[3]} - \alpha * \dfrac{\partial \mathcal{L}}{\partial W^{[3]}}$
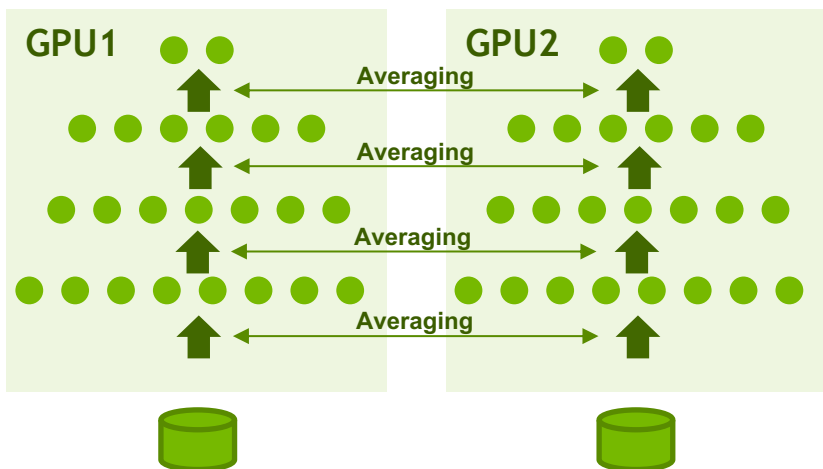
$W^{[2]}$

$W^{[2]} = W^{[2]} - \alpha * \dfrac{\partial \mathcal{L}}{\partial W^{[2]}}$

$W^{[1]}$

$W^{[1]} = W^{[1]} - \alpha * \dfrac{\partial \mathcal{L}}{\partial W^{[1]}}$

CPU/GPU

sgd
momentum
nag
adagrad
adadelta
rmsprop

1. Read the data
2. Transport the data
3. Pre-process the data
4. Queue the data
5. Transport the data
6. Calculate activations for layer one
7. Calculate activations for layer two
8. Calculate the output
9. Calculate the loss
10. Backpropagate through layer three
11. Backpropagate through layer two
12. Backpropagate through layer one
13. Execute optimization step
14. Update the weights
15. Return control

NVIDIA. DEEP LEARNING INSTITUTE

# TRAINING A NEURAL NETWORK

## Multiple GPUs



GPU

$\hat{y}$

$\mathcal{L}(\hat{y}, y)$

$W^{[3]}$

$W^{[3]} = W^{[3]} - \alpha * \dfrac{\partial \mathcal{L}}{\partial W^{[3]}}$

$W^{[2]}$

$W^{[2]} = W^{[2]} - \alpha * \dfrac{\partial \mathcal{L}}{\partial W^{[2]}}$

$W^{[1]}$

$W^{[1]} = W^{[1]} - \alpha * \dfrac{\partial \mathcal{L}}{\partial W^{[1]}}$

CPU/GPU

CPU/GPU

GPU

$\mathcal{L}(\hat{y}, y)$

$\hat{y}$

$W^{[3]} = W^{[3]} - \alpha * \dfrac{\partial \mathcal{L}}{\partial W^{[3]}}$

$W^{[3]}$

$W^{[2]} = W^{[2]} - \alpha * \dfrac{\partial \mathcal{L}}{\partial W^{[2]}}$

$W^{[2]}$

$W^{[1]} = W^{[1]} - \alpha * \dfrac{\partial \mathcal{L}}{\partial W^{[1]}}$

$W^{[1]}$

CPU/GPU
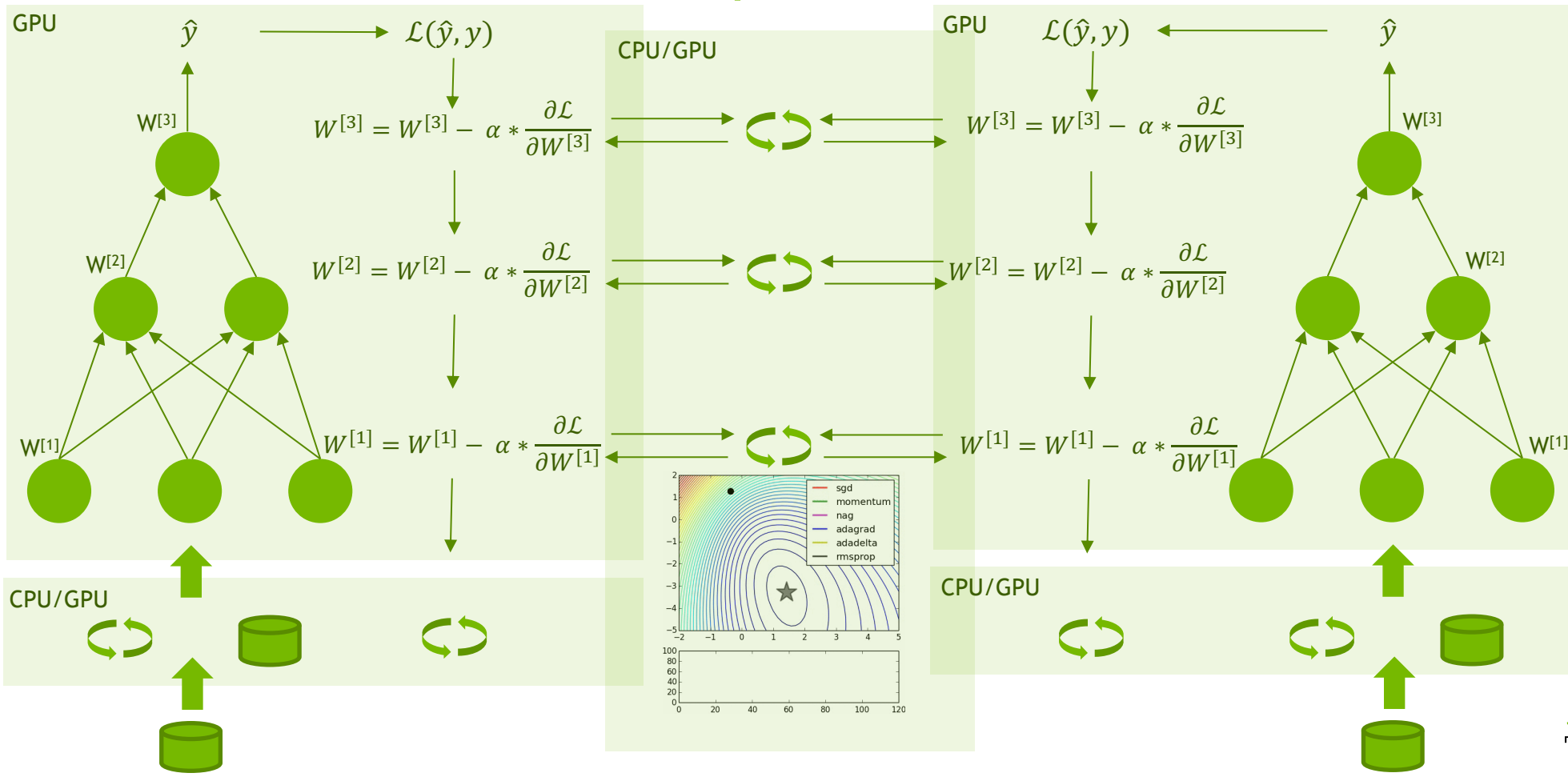
# FUNDAMENTALS OF DEEP LEARNING FOR MULTI-GPUS

## LAB 2, PART 1: INTRODUCTION TO HOROVOD

**NVIDIA** | DEEP LEARNING INSTITUTE

# TRAINING A NEURAL NETWORK
## Multiple GPUs

# MEET HOROVOD

Library for distributed DL

Works with stock TensorFlow, Keras, PyTorch, and MXNet

Installs with pip

Uses advanced algorithms; leverages high-performance networks (RDMA, GPUDirect).



36

horovod.ai

# MEET HOROVOD

Infrastructure team provides container and MPI environment

ML engineers use DL frameworks that they love

Both have consistent expectations for distributed training across frameworks

HOROVOD

horovod.ai

# USING HOROVOD

# INITIALIZE THE LIBRARY
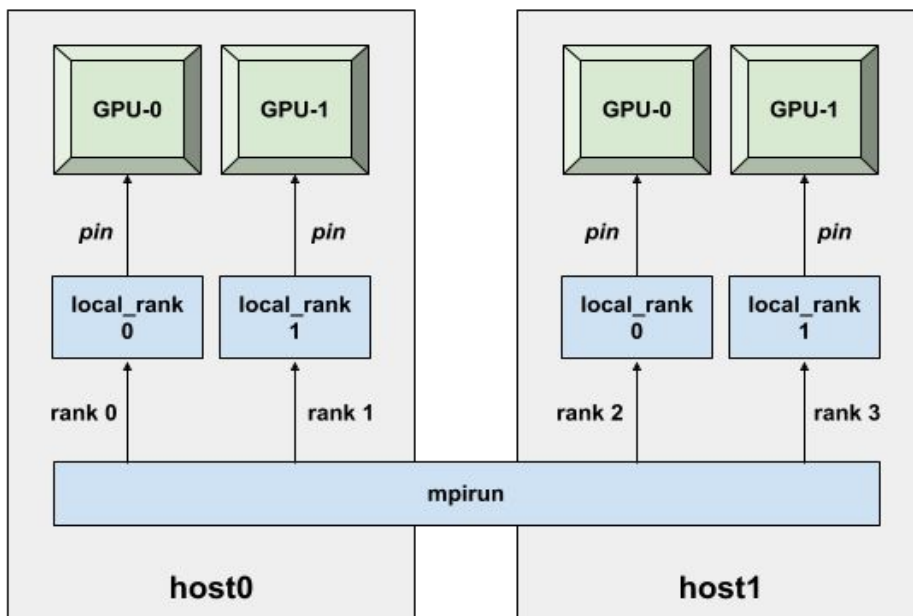
```
import horovod.tensorflow.keras as hvd

hvd.init()
```

39

# PIN GPU TO BE USED

```
gpus = tf.config.experimental.list_physical_devices('GPU')
if gpus:
    tf.config.experimental.set_memory_growth(gpus[hvd.local_rank()], True)
    tf.config.experimental.set_visible_devices(gpus[hvd.local_rank()], 'GPU')
```



40

# ADD DISTRIBUTED OPTIMIZER

```
opt = hvd.DistributedOptimizer(opt)
```
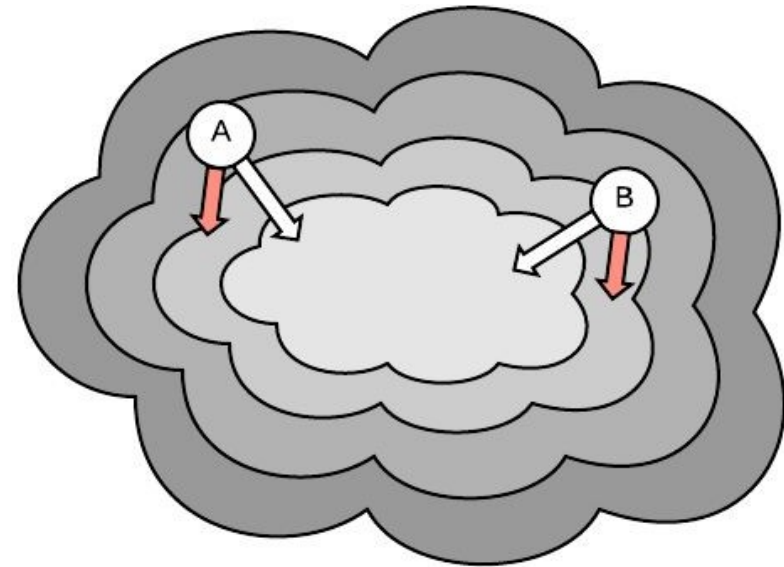
# SYNCHRONIZE INITIAL STATE

```
callbacks.append(hvd.BroadcastGlobalVariablesCallback(0))
```

```
model.fit(..., callbacks, ...):

        ...
```

# CHECKPOINT ONLY ON ONE WORKER

```
checkpoint_callback = tf.keras.callbacks.ModelCheckpoint(...)


if hvd.rank() == 0:

        callbacks.append(checkpoint_callback)



model.fit(..., callbacks, ...):
```

```
        ...
```

# DATA PARTITIONING: OPTION 1

Shuffle the dataset

Partition records among workers

Train by sequentially reading the partition

After epoch is done, reshuffle and partition again

Dataset

Shuffle

Dataset

Worker

Worker

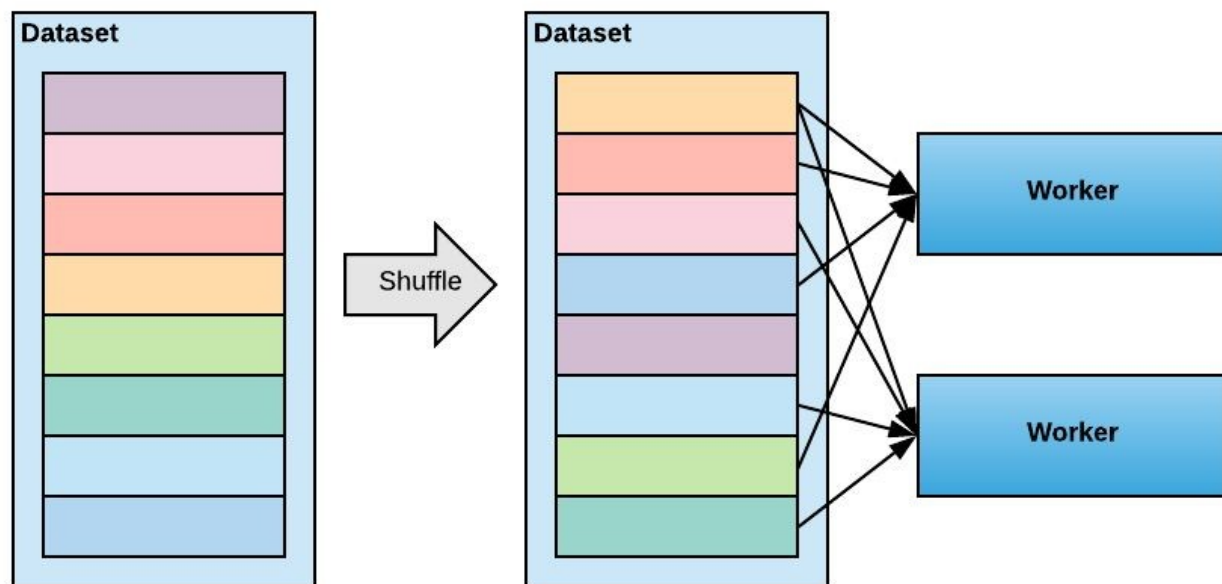**NOTE:** make sure that all partitions contain the same number of batches, otherwise the training will deadlock

# DATA PARTITIONING: OPTION 2

Shuffle the dataset

Train by randomly reading data from whole dataset

After epoch is done, reshuffle

# HOROVOD FOR ALL

```
import horovod.tensorflow as hvd
import horovod.tensorflow.keras as hvd
import horovod.torch as hvd
import horovod.mxnet as hvd
```

# RUNNING HOROVOD

## Single-node:

```
$ mpirun -np 4 python train.py
```

## Multi-node:

```
$ mpirun -np 8 -H server1:4,server2:4 python train.py
```

# FUNDAMENTALS OF DEEP LEARNING FOR MULTI-GPUS

## LAB 3, PART 1: SCALING THE BATCH SIZE

NVIDIA | DEEP LEARNING INSTITUTE

# CAN WE INCREASE THE BATCH SIZE INDEFINITELY?

# IN TERMS OF IMAGES / SECOND?

## Yes



(a) Tiramisu

(b) DeepLabv3+

Kurth, T., Treichler, S., Romero, J., Mudigonda, M., Luehr, N., Phillips, E., ... & Houston, M. (2018, November). Exascale deep learning for climate analytics. In Proceedings of the International Conference for High Performance Computing, Networking, Storage, and Analysis (p. 51). IEEE Press. arXiv:1810.01993

# IN TERMS OF STEPS TO CONVERGENCE?

## There are limits



(a) Simple CNN on MNIST
(b) Simple CNN on Fashion MNIST
(c) ResNet-8 on CIFAR-10
(g) Transformer on Common Crawl

(d) ResNet-50 on ImageNet
(e) ResNet-50 on Open Images
(f) Transformer on LM1B
(h) VGG-11 on ImageNet

(i) LSTM on LM1B

Shallue, C. J., Lee, J., Antognini, J., Sohl-Dickstein, J., Frostig, R., & Dahl, G. E. (2018). Measuring the effects of data parallelism on neural network training. arXiv:1811.03600

# IN TERMS OF STEPS TO CONVERGENCE?

## There are limits

# IMPACT ON ACCURACY

## Naïve approaches lead to degraded accuracy



ImageNet by ResNet50 without Data Augmentation

- Batch=16k, LARS
- Batch=16k, no LARS
- Batch=256

ImageNet by ResNet50 without Data Augmentation

- Batch=32k, LARS
- Batch=32k, no LARS
- Batch=256

You, Y., Zhang, Z., Hsieh, C., Demmel, J., & Keutzer, K. (2017). ImageNet training in minutes. arXiv: 1709.05011

# IMPACT ON ACCURACY

## Naïve approaches lead to degraded accuracy



(a) Training error

(b) Validation error

Hoffer, E., Hubara, I., & Soudry, D. (2017). Train longer, generalize better: closing the
generalization gap in large batch training of neural networks. arXiv:1705.08741

# IMPACT ON ACCURACY
## Why? Generalization and flatness of minima?



Training Function

Testing Function

$f(x)$

Flat Minimum

Sharp Minimum

Keskar, N. S., et al. (2016). On large-batch training for deep learning: Generalization gap and sharp minima. arXiv:1609.04836

# IMPACT ON ACCURACY

## Why does it happen? Noise in the gradient update.



Keskar, N. S., et al. (2016). On large-batch training for deep learning: Generalization gap and sharp minima. arXiv:1609.04836

# IMPACT ON ACCURACY



(a) 0.0, 128, 7.37%   (b) 0.0, 8192, 11.07%   (c) 5e-4, 128, 6.00%   (d) 5e-4, 8192, 10.19%

(e) 0.0, 128, 7.37%   (f) 0.0, 8192, 11.07%   (g) 5e-4, 128, 6.00%   (h) 5e-4, 8192, 10.19%

Figure 3: The 1D and 2D visualization of solutions obtained using SGD with different weight decay and batch size. The title of each subfigure contains the weight decay, batch size, and test error.

Li, H., Xu, Z., Taylor, G., & Goldstein, T. (2017). Visualizing the Loss Landscape of Neural Nets. arXiv:1712.09913

# FUNDAMENTALS OF DEEP LEARNING FOR MULTI-GPUS

## LAB 3, PART 2: OPTIMIZATION STRATEGIES

# WHAT CAN WE DO TO IMPROVE THE OPTIMIZATION PROCESS?

- Manipulate the learning rate?

- Add noise to the gradient?

- Manipulate the batch size?

- Change the learning algorithm?

# WHAT CAN WE DO ABOUT IT?
## Early approaches: scaling the learning rate

"Theory suggests that when multiplying the batch size by k, one should multiply the learning rate by $\sqrt{}$(k) to keep the variance in the gradient expectation constant.

$$\mathrm{cov}\left(\Delta\mathbf{w}, \Delta\mathbf{w}\right) \approx \frac{\eta^2}{M}\left(\frac{1}{N}\sum_{n=1}^{N}\mathbf{g}_n\mathbf{g}_n^\top\right) \longrightarrow \eta \propto \sqrt{M}$$

...

Theory aside, for the batch sizes considered in this note, the heuristic that I found to work the best was to multiply the learning rate by k when multiplying the batch size by k. I can't explain this discrepancy between theory and practice."

In practice linear scaling is still frequently used.

Krizhevsky, A. (2014). One weird trick for parallelizing convolutional neural networks. arXiv:1404.5997

# WHAT CAN WE DO ABOUT IT?
## Warmup strategies

- A lot of networks will diverge early in the learning process

- Warmup strategies address this challenge

**Gradual warmup.** We present an alternative warmup that *gradually* ramps up the learning rate from a small to a large value. This ramp avoids a sudden increase of the learning rate, allowing healthy convergence at the start of training. In practice, with a large minibatch of size $kn$, we start from a learning rate of $\eta$ and increment it by a constant amount at each iteration such that it reaches $\hat{\eta} = k\eta$ after 5 epochs (results are robust to the exact duration of warmup). After the warmup, we go back to the original learning rate schedule.

Goyal, P., Dollár, P., Girshick, R., Noordhuis, P., Wesolowski, L., Kyrola, A., ... & He, K. (2017). Accurate, Large Minibatch SGD: Training ImageNet in 1 Hour. arXiv:1706.02677

# WHAT CAN WE DO ABOUT IT?

## Batch Normalization

Batch normalization improves the learning process by minimizing drift in the distribution of inputs to a layer

It allows higher learning rates and reduces the need to use dropout

The idea is to normalize the inputs to all layers in every batch (this is more sophisticated than simply normalizing the input dataset)



Figure 1: (a) *The test accuracy of the MNIST network trained with and without Batch Normalization, vs. the number of training steps. Batch Normalization helps the network train faster and achieve higher accuracy.* (b, c) *The evolution of input distributions to a typical sigmoid, over the course of training, shown as {15, 50, 85}th percentiles. Batch Normalization makes the distribution more stable and reduces the internal covariate shift.*

Ioffe and Szegedy (2015). Batch Normalization: Accelerating Deep Network Training by Reducing Internal Covariate Shift. arXiv:1502.03167

DEEP LEARNING INSTITUTE

# WHAT CAN WE DO ABOUT IT?

## Ghost Batch Normalization

- The original batch normalization paper suggests using the statistics for the entire batch, but what should that mean when we have multiple GPUs?

- We can introduce additional noise by calculating smaller batch statistics ("ghost batches").

- Batch normalization is thus carried out in isolation on a per-GPU basis.

Hoffer, E., Hubara, I., & Soudry, D. (2017). Train longer, generalize better: closing the generalization gap in large batch training of neural networks. arXiv:1705.08741

# WHAT CAN WE DO ABOUT IT?
## Adding noise to the gradient

- Keeps the covariance constant with changing batch size (as $\sigma^2 \propto M$)
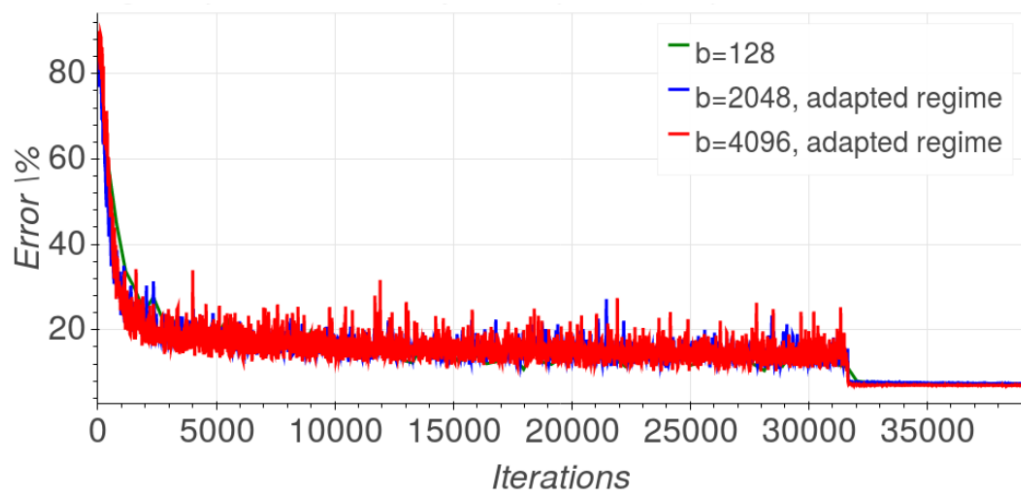
- Does not change the mean

Furthermore, we can match both the first and second order statistics by adding multiplicative noise to the gradient estimate as follows:

$$\hat{\mathbf{g}} = \frac{1}{M} \sum_{n \in B}^{N} \mathbf{g}_n z_n \,,$$

where $z_n \sim \mathcal{N}\left(1, \sigma^2\right)$ are independent random Gaussian variables for which $\sigma^2 \propto M$. This can be verified by using similar calculation as in appendix section A. This method keeps the covariance constant when we change the batch size, yet does not change the mean steps $\mathbb{E}\left[\Delta \mathbf{w}\right]$.

Hoffer, E., Hubara, I., & Soudry, D. (2017). Train longer, generalize better: closing the generalization gap in large batch training of neural networks. arXiv:1705.08741

# WHAT CAN WE DO ABOUT IT?

## Longer training with larger learning rate



(a) Validation error

(b) Validation error - zoomed

Hoffer, E., Hubara, I., & Soudry, D. (2017). Train longer, generalize better: closing the generalization gap in large batch training of neural networks. arXiv:1705.08741

# WHAT CAN WE DO ABOUT IT?

## Increasing the batch size, instead of learning rate decay



(a)

(b)

Smith, S. L., Kindermans, P. J., & Le, Q. V. (2017). Don't Decay the Learning Rate, Increase the Batch Size. arXiv:1711.00489

# WHAT CAN WE DO ABOUT IT?

## LARS: Layer-wise Adaptive Rate Scaling



Figure 2: LARS: local LR for different layers and batch sizes

You, Y., Gitman, I., & Ginsburg, B. Large batch training of convolutional networks. arXiv:1708.03888

# WHAT CAN WE DO ABOUT IT?
## LARS: Layer-wise Adaptive Rate Scaling

Control magnitude of the layer $k$ update through local learning rate $\lambda_k$:

$$\Delta w_k(t + 1) = \lambda_k * G_k(w(t))$$

where:

$G_k(w(t))$: stochastic gradient of $L$ with respect to $w_k$,

$\lambda_k$: local learning rate for layer $k$, defined as

$$\lambda_k = min(\gamma, \ \eta \cdot \frac{||w_k(t)||_2}{||G_k(w(t))||_2})$$

where

$\eta$ is trust coefficient (how much we trust stochastic gradient)

$\gamma$ is global learning rate policy (steps, exponential decay, ...)

You, Y., Gitman, I., & Ginsburg, B. Large batch training of convolutional networks. arXiv:1708.03888

# WHAT CAN WE DO ABOUT IT?

LARC: Layer-wise learning rates with clipping; SGD with momentum is base optimizer

LAMB: Layer-wise learning rates; Adam as base optimizer

- More successful than LARC at language models like BERT

NovoGrad: Moving averages calculated on a per-layer basis

- Also useful in several different domains

# FUNDAMENTALS OF DEEP LEARNING FOR MULTI-GPUS

## LAB 3: INTRODUCTION TO THE ASSESSMENT

# A FINAL ASSESSMENT TO TEST YOUR SKILLS

This assessment will test all of what you have learned in this course. You are required to take a serial training script, convert it to use Horovod, and obtain a target training and validation accuracy in a fixed amount of time.

The training is very similar, but this time we are using CIFAR-10 instead of Fashion MNIST.

# A FINAL ASSESSMENT TO TEST YOUR SKILLS

You can make changes to the assessment.py script in the JupyterLab environment and test the performance in the notebook.

When you are done, go back to the browser tab you launched this lab from, and click "Assess".

You will get output after a few minutes indicating whether you passed. If not, go back and try again! Good luck ☺

# FINAL THOUGHTS

Use NGC containers for high-performance, multi-GPU training.



### Innovate Faster

Get up and running quickly while reducing the complexity typically associated with setting up software.
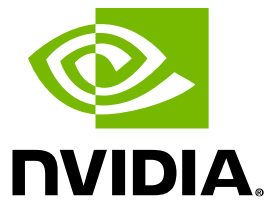


### Stay Up to Date

The top deep learning containers are updated monthly to keep your systems running at peak performance. All containers provide easy access to fully-tested and optimized software releases.



### Run Anywhere

NGC containers are built to run on-prem, in the cloud, or in hybrid deployments with Docker and Singularity runtimes. This allows for maximum utilization of available GPUs, portability, and scalability.

Please take the survey!